

# HYPertext AND THE OXFORD ENGLISH DICTIONARY

*Hypertext databases can be produced by converting existing text documents to electronic form. The basic task in conversion is identification of fragments. We illustrate that this is not always a straightforward process with an analysis of the Oxford English Dictionary.*

DARRELL R. RAYMOND *and* FRANK WM. TOMPA

Low-cost high-capacity storage devices are making widespread computerized access to existing reference texts a reality. Documents currently available in electronic form include Bowker's *Books In Print*, Microsoft's *Bookshelf* (a compendium of ten works including *Roget's Thesaurus*, *Bartlett's Familiar Quotations*, the *World Almanac*, and the *U.S. Zip Code Directory*), the *McGraw-Hill Concise Encyclopedia of Science and Technology*, and the *Electronic Holmes Companion*. Academics can make use of the approximately one thousand texts maintained by the Oxford Text Archives, which include documents of scholarly interest such as the *Tibetan Book of the Dead* and *Steppenwolf*. The basic assets of computerization—speed of access and quantity of storage—make almost any form of electronic documentation useful. In addition, computerization presents new opportunities in advanced representations of such documents, for both storage and display. Hypertext is one such opportunity.

Hypertext is a database technology typically characterized by small fragments of text interconnected by machine-supported links [5, 10]. Since most existing

texts are not fragmented in this way, the key problem in converting them to hypertext is the development of text fragments and links. Fragmentation must be carried out carefully if the hypertext is to be a faithful representation of the original. If the document was created decades or centuries ago, its authors or editors may not be available for consultation. As a result, structure must be inferred from a careful study of the text and from consultation with existing experts. The characteristics of pen and paper or the printing press often exert a powerful influence on a document that should be controlled if possible. Finally, there is always the danger that implicit, unrecognized structure will be lost when converting to a new representation.

This article reports our experiences in the evaluation of a potential hypertext representation of the *Oxford English Dictionary*. With the completion of the *Supplement* to the *OED*, the Oxford University Press has taken under consideration an electronic form for the dictionary [18]. The text of the *OED* has been keyboarded by an independent firm, and we are currently designing suitable computer-based tools for storage, editing, and searching [3]. Our work on the *OED* provides some insight into the problems of converting existing documents to hypertext.

## THE OXFORD ENGLISH DICTIONARY

The *Oxford English Dictionary* [13] is the largest and most scholarly dictionary of written English. Its production spanned the period from 1884 to 1928, with nearly 30 years of preliminary effort in planning and collecting material. In its standard form, the *OED* consists of 12 books containing 41.81 million words in 252,259 entries, and 1.86 million quotations. An important subsequent work is the four-volume *Supplement* which was produced from 1958 to 1986. The *Supplement* contains 69,372 entries, 14.5 million words, and 560,000 quotations. Seventy percent of the *Supplement* entries describe words not covered in the *OED*; the remaining 30 percent contain emendations and additions for existing *OED* entries.

The basic unit of the *OED* is the entry; each entry details the historical development of a given word. A typical *OED* entry is shown for the word **abbreviate**, in Figure 1.

An entry consists of a main form (the word being defined), its pronunciation, its part of speech, a list of variant forms, the etymology, and a list of senses. Each sense contains a definition and usually some illustrative quotations chosen from different sources and time periods. Senses often have nested subsenses (e.g., those labelled **b.** and **c.** in Figure 1) in the case of commonly used words, prefixes, or suffixes.

**Abbreviate** (ăbrĭ-vi,ăt), *v.*, also 5-7 **abreviate**. [f. ABBREVIATE *ppl. a.*; or on the analogy of *vbs.* so formed; see -ATE. A direct representative of *L. abbreviāre*; as ABRIDGE, and the obs. ABREVI, represent it indirectly, through OFr. *abregier* and mid. Fr. *abrévier*. Like the latter, *abbreviate*, was often spelt *a-breviate* in 5-7.] To make shorter, shorten, cut short in any way.  
 1530 PALSGR., I *abrevyate*: I make a thynge shorte, *Je abregre*.  
 1625 BACON *Essays* xxiv. 99 (1862) But it is one Thing to Abbreviate by Contracting, Another by Cutting off.  
 † **1.** *trans.* To make a discourse shorter by omitting details and preserving the substance; to abridge, condense. *Obs.*  
 c 1450 *Chester Pl. I. 2* (Sh. Soc.) This matter he abbreviated into plays twenty-foure. 1592 GREENE *Comy catching* III. 16 The queane abreviated her discourse. 1637 RALEIGH *Mahomet* 34 Abbreviated out of two Arabique writers translated into Spanish. 1672 MANLEY *Interpreter* pref., I have omitted several Matters .. contracted and abbreviated Others.  
 † **b.** To make an abstract or brief of, to epitomize. *Obs.*  
 c 1450 TREVISA *Higden's Polychr. I. 21* (Rolls Ser.) Trogus Pompeius, in hys xli<sup>th</sup> iiii. bookes, allemoste of alle the storyes of the worlde, whom Iustinus his disciple did abbreviate. 1603 FLORIO *Montaigne* (1634) 627 To reade, to note, and to abbreviate Polibius. 1648-9 *The Kingdomes Weekly Intelligencer* Jan. 16 to 23 The high court of Justice did this day sit again concerning the trial of the King. The charge was brought in and abbreviated.  
 † **c.** *Math.* To reduce (a fraction) to lower terms. *Obs.*  
 1796 *Mathem. Dict. I. 2* To abbreviate fractions in arithmetic and algebra, is to lessen proportionally their terms, or the numerator and denominator.  
 † **2.** *intr.* To speak or write briefly, to be brief. *Obs.*

FIGURE 1. Entry for Abbreviate

While some of the elements in an entry are indicated by special symbols (e.g., etymologies are surrounded by square brackets), most are indicated only implicitly by position and font. In order to preserve its structure, the text of the *OED* was manually keyboarded, with insertion of typographical tags to capture the most salient characteristics of the entries. These tags were designed for ease of data entry rather than completeness of structure representation, and so extensive work was undertaken to build a parser that could complete the tagging and verify the data [11]. The tagged and parsed text for *abbreviate*, up to and including the first two quotations, is shown in Figure 2. These tags are currently being exploited in the development of a database index and are crucial in identifying potential components of a hypertext.

## WHY HYPERTEXT FOR THE OED?

The main reason for considering a hypertext representation of the *OED* is to support browsing. The *OED* can be treated as a text database to which formal queries are posed, (e.g., *What interjections were in common use in the period 1670-1720?*) [2, 7]. Nevertheless, experience with an extremely rapid search program called PAT [8] has shown that browsing through the dictionary by means of fixed string searches is an invaluable adjunct to formal querying, and is often more fruitful, serendipitous, and enjoyable. Browsing is a two-stage process: users specify a pattern to be located and then navigate in the vicinity of the pattern. Our users have typically employed PAT to search for interesting phrases or words, and then consulted the relevant entries in the original paper dictionary. This switch of medium is necessary because PAT is only somewhat knowledgeable about the structure of the dictionary, and does not support a sufficient navigational browsing capability. An obvious extension is a hypertext-like browsing facility as a front end to PAT.

A second reason for hypertext representation is to explore alternative methods for entry display. The *OED* was originally contracted for a specific number of pages, which the editors soon learned would be insufficient. As a result, they employed every means at their disposal to control the size of the *OED*, short of reducing its quality or comprehensiveness [14]. The need to maximize the use of space led to very dense typesetting and the extensive use of abbreviations and symbols. Spatial techniques were restricted to paragraphs and different sizes of typeface. While the result is admirable, we need not observe this constraint any longer. A hypertext representation should provide better visual salience and more rapid navigation around large or alphabetically-distant entries. A key mechanism to be employed is dynamic reformatting of entries according to user specification.

A third reason for hypertext representation is to integrate the *OED* with the users' tasks. Our users typically query and browse the *OED* as part of more extended tasks. At a minimum, users want to save their results and queries for use in future sessions, but they also

```

<entry> <hwgp> <hwlem>abbreviate </hwlem> <pron id=0000041884>a&breve.br
<i>i&mac.</i> <sd.vi&syllab.<i>e</i> <su>i</su> t </pron>, <pos>v.</pos> </hwgp>
<vfl> Also <vd>5&en.7</vd> <vf>abbreviate</vf>. </vfl> <etym> f.<xra
id=0000041880><xlem>abbreviate</xlem> <pos>ppl. a.</pos> </xra>; or on the
analogy of vbs. so formed; see <xra id=0000041881> <xlem>-ate</xlem> </xra>.
&es.A direct representative of L. <cf>abbrevia&mac.re</cf>; as <xra
id=0000041882><xlem>abridge</xlem> </xra>, and the obs. <xra
id=0000041883><xlem>abrey</xlem> </xra>, represent it indirectly, through OFr.
<cf>abregier</cf> and mid.Fr. <cf>abre&acu.vier</cf>. &es.Like the latter,
<cf>abbreviate </cf>, was often spelt <cf>a-breviate</cf> in 5&en.7. </etym>
<sen4> <sen6> To make shorter, shorten, cut short in any way. <qpara> <quot>
<qdat>1530</qdat> <auth>Palsgr.</auth>, <txt>I abreyate: I make a thyng
shorte, <i>Je abreye</i>. </txt> </quot> <quot> <qdat>1625</qdat>
<auth>Bacon</auth> <wk>Essays</wk> xxiv. 99 (1862) <txt>But it is one thing to
Abbreviate by Contracting, Another by Cutting off.</txt> </quot> </qpara>
</sen6> </sen4>

```

FIGURE 2. Tagged Data for Abbreviate

expect access to annotation facilities, the ability to cut and paste fragments of *OED* text into other documents, routines for sorting and filtering extracted quotations, and tools for statistical analysis of selected variables. Hypertext can facilitate a consistent and simple interface to this wide range of tools.

### CONVERTING TEXT TO HYPERTEXT

Hypertext's main characteristic is commonly held to be its nonlinear network of interconnections. Contemplating the conversion of the *OED* to hypertext has led us to address the problem of defining nodes and links. A document can be broken into many arbitrary networks of nodes, but few of these are suitable representations of the original text. The selection of nodes and links has a significant impact on the usability of the hypertext and its validity as a representation of the original document.

From the point of view of document conversion, hypertext's main characteristic is *fragmentation*. Fragments are pieces of content or structure which are both discrete and independent. By discrete we mean that the distinction between components is explicit and well-defined; by independent we mean that components are capable of standing in multiple relations to one another. Ideally, a hypertext node is a text fragment which contains a single, independent concept; common examples are footnotes, references, and annotations. Similarly, a hypertext link is a relationship with a specific source and destination node, often made concrete with an arrow and a label. Usually a link must be explicitly selected by the user to be traversed.

Texts often contain ideas that are not confined to discrete fragments. For example, Hamlet's indecisiveness is a theme based on a set of actions, some explicit, some only alluded to, and some conspicuous by their absence. Making these actions explicit is a highly subjective process; moreover, it involves a radical departure from the nature of the text. Hypertext proponents sometimes insist that most texts would be better off as hypertext, and that linear structure is primarily a limi-

tation of the media. We suggest that linear structure is a virtue achieved at some cost by the document's creator. Furthermore, a subtle treatment of a theme may communicate an idea more artistically (and hence more successfully) than any explicit statement.

Thus for some purposes and documents, hypertext representations are inadequate. Fragmentation is intended to make an implicit structure explicit, so the key question in conversion must be: will an explicit structure be as expressive as the implicit structure? When the answer is yes, the document will gain from conversion; otherwise, conversion will degrade the representation of the document.

### CHARACTERISTICS OF THE *OED*

Much of the *OED* consists of quotations that have been extracted from other texts. These quotations exhibit sufficient independence to be detached from their source, yet are still illustrative of a particular sense of a word. Fragmentation is not a problem if only a quotation database is considered. Complete entries are also independent, since they are generally intended to be read as individual entities. Searches, however, are not often satisfied by exactly one entry. Some entries are extremely large, many have important cross-references, and often homonyms must be examined for related explanations.

The subcomponents of the entries (e.g., senses, etymologies) adhere to a formal structure, since we can describe them with a context-free grammar. The Oxford editors have taught us, however, that an entry is not merely a collection of independent subcomponents [15]. Entries are stylistic, creative wholes in which the relative sizes and arrangement of the parts often convey important (although implicit) information about the development of a word. The fragmentation of an entry may camouflage or destroy this implicit structure.

An attempt to represent each entry as a discrete fragment is complicated by the variance in entry size and structure. Entries in the *OED* range from "Gig: see Jig" to the entry for *set v.*, which requires almost half a

megabyte of storage. Some entries have hundreds of quotations, but many have none. An entry may contain a well-balanced arrangement of senses and quotations, or it may consist largely of possible compounds or formations, (e.g., the entry for *un-*). While some entries can be read quickly, viewing even the structural skeleton of long entries can be time-consuming. Clearly, the attempt to find a representation for such a wide range will not be a simple task.

The range is wide, but what is the distribution? A histogram of the frequency of entry size is shown in Figure 3. This figure displays the frequency of entries of size less than 5,000 bytes (entry size has been rounded to the nearest ten characters).

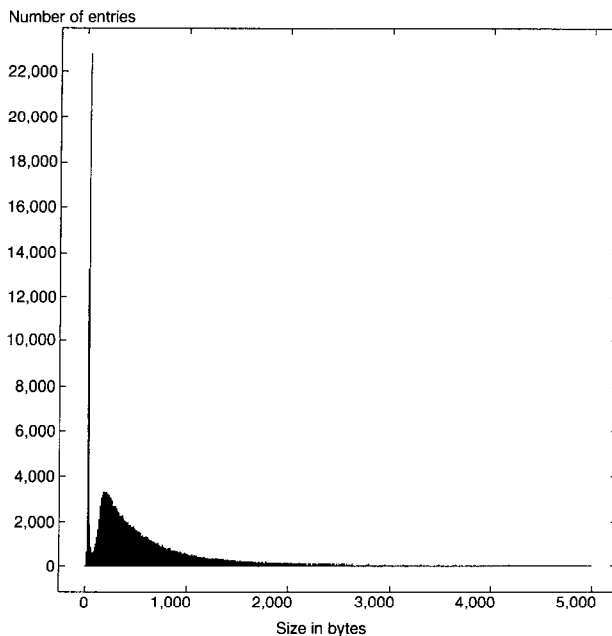


FIGURE 3. Distribution of Entry Size

The most striking feature of the distribution is that approximately 50,000 entries (one-fifth of the total) contain fewer than 50 characters. The majority of these are entries for infrequently used variants of words and cross-references to other entries. Ninety-five percent of all entries are smaller than 4000 characters, with 20 percent of all entries larger than the mean of 1060 characters. It should be noted that some entries will become larger when the material in the *Supplement* is integrated with that in the *OED*.

It would appear that large entries are not extremely common, and so the problem of representation may be less severe than it first appeared. The probability of access, however, is not equal for every entry. The size of an entry is largely dependent on the number of suitable quotations submitted. It should not be surprising that many of the smaller entries are for obsolete or infrequently used words, while commonly used words

have larger and more complicated entries. If the commonly used words are also more likely to be browsed, the need for this kind of representation increases. In fact, an argument could be made that the larger entries are more likely to be browsed whether they are common or not, since their structure and content are more difficult to perceive and remember.

The probability of accessing large entries is also increased by the nature of our searching tools. For many search patterns there are too many matches to be returned to the user; in these cases PAT can be instructed to give a random sample of the results. The probability of selecting a large entry in a random probe is proportional to the fraction of the *OED* covered by large entries. Only five percent of entries are larger than 4000 characters, but they account for 48 percent of the bulk of the *OED*. The probability of matching to a large entry in a random sample of ten matches is  $1 - (1 - 0.48)^{10}$ , a virtual certainty. Empirical observations confirm this analysis; samples of matches nearly always contain one or more large entries.

The display of larger entries might be facilitated by structural views or abbreviations such as those described by Furnas [6]. Structural information can be extracted from the tags and employed in the construction of a structural view. Figure 4 shows the relationship between the number of tags and the size of entries. The graph shows generally that the larger the entry, the more tags it contains.

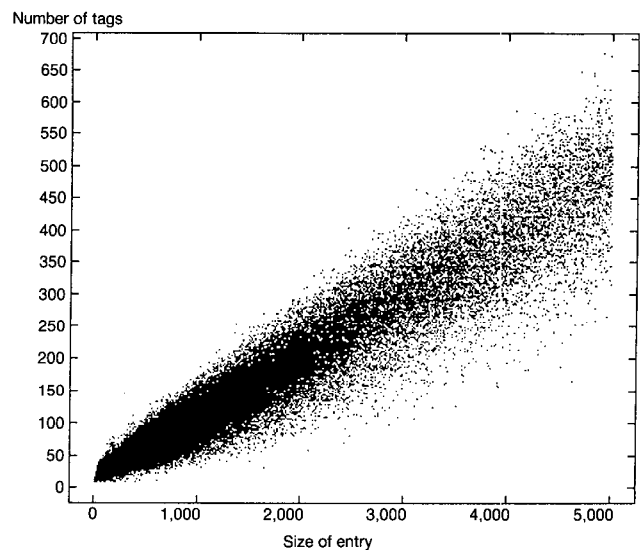


FIGURE 4. Distribution of Tags

Many tags, however, enclose short pieces of text (e.g., a single part of speech, a variant form, or a quotation date). Such tags designate a flat structure that does not facilitate abbreviated display. Since the sense tags constitute the bulk of the nested structural information, restricting our attention to them gives a better appreciation of the substructure. This data is shown in Figure 5.

Clearly the number of sense tags is not well correlated with entry size. Some entries have a large number of quotations or long etymological notes but a very shallow sense structure; a sense-oriented view of such an entry is not indicative of its size or comprehensiveness. On the other hand, some entries have a deeper hierarchy with relatively little leaf content; in this case a sense-oriented skeletal view overemphasizes the content of the entry.

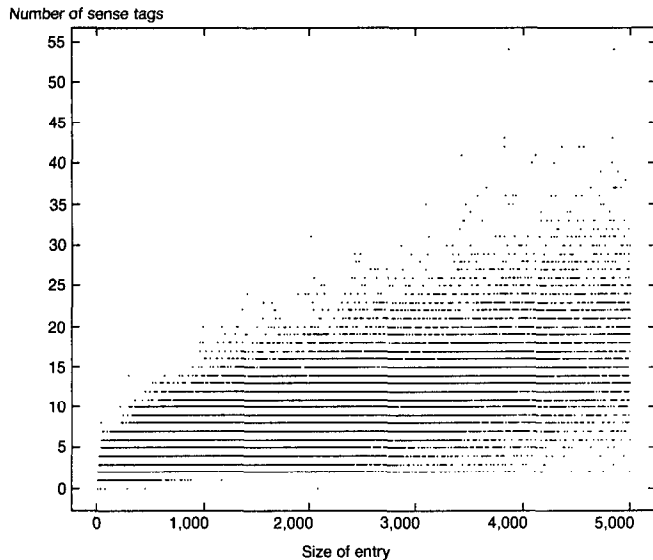


FIGURE 5. Distribution of Sense Tags

From the previous analysis it seems that static fragmentation will not be a satisfactory solution, both because of the possibility of losing implicit information and because no static fragmentation seems appropriate for all entries. As a result, we have chosen to support dynamic fragmentation, controlled by the user. By permitting various types of elision and formatting in a rapid response display tool, users can employ multiple fragmentations while viewing an entry. The results of a prototype formatter are shown in Figure 6. Here "abbreviate" is displayed as *OED*-style proof, a sense structure skeleton, the sense content without quotations, and quotations only.

Turning to the consideration of links, first we must ascertain the characteristics of the links that are already present in the data and tags. The most explicit links are the cross-references, which are pointers to other entries. These can appear in etymological notes or sense text, and are visually designated by printing the main form of the cross-referenced entry in small capitals. In Figure 1 there are cross-references to *-ATE*, *ABRIDGE*, and *ABREVVY*. The *OED* contains 569,000 cross-references for an average of 2.26 per entry, making these a substantial source of potential hypertext links.

It is interesting to investigate the distribution of cross-reference destinations. Figure 7 contains four

graphs showing the distribution of cross-references for each letter of the alphabet, normalized according to the number of entries per letter. Figure 7a shows cross-references to previous letters, Figure 7b shows cross-references to the same letter, Figure 7c shows cross-references to later letters, and Figure 7d shows cross-references to suffixes. Most cross-references are to words beginning with the same letter of the alphabet, with the major exception being cross-references to suffixes. As might be expected, the number of cross-references to previous letters increases through the alphabet, while the number of cross-references to later letters decreases. It should be noted, however, that there are several rather odd exceptions to the general trends. These exceptions may be due to the influence of the editor, the lexicographic position of the letter, or etymological characteristics common to words beginning with that letter.

The general trend of the distribution of cross-references is a result of several factors. First, cross-references in the *OED* usually point to words with similar spelling. Many cross-references in definitional text are prefaced with *erroneous spelling of*, *variant of*, *obs. form of*, *verbal form of* and other such phrases that typically indicate a word with very similar spelling. The major exception to this rule is suffixes, which are referenced more evenly from the whole alphabetic range. Cross-references in etymologies are often of the form *root + suffix*, for example *Musal*. . [f. *MUSE sb.*<sup>1</sup> + *-AL*]. A second factor is the history of the development of the dictionary. Volumes were generally compiled in alphabetical order, so the compilers of the later volumes had more information on which to base cross-references, and were more likely to cross-reference existing entries than the still uncompiled ones. A third factor is psychological. Since an editor would be most aware of the section of the *OED* currently in progress, he or she would be more likely to insert cross-references to entries in that section.

A more general type of link is known as a *lexical link*. These links result from the simple observation that since the *OED* is a book that defines most English words, every word in the *OED* can be seen as the source of a possible link to its definition. The *OED* and other comprehensive dictionaries are thus unique in containing to some extent their own reference material. *Lexical links need not be confined to the OED*; it would be highly desirable to link words in other texts to their *OED* entry. Nevertheless, automatically determining the source and destination of lexical links is not simple. For example, consider the definition in sense 1. of *Fossic v.*:

To search for gold by digging out crevices with knife and pick, or by working in washing-places and abandoned workings in the hope of finding particles or small nuggets overlooked by others. Also, to *fossick about*.

The first problem is to determine which lexical items should be linked to their definitions. For instance,

abbreviate v.

**abbreviate**, v. Also 5-7  
 abbreviate. [f. abbreviate *ppl. a.*; or on the analogy of vbs. so formed; see -ate. A direct representative of L. *abbreviare*; as *abridge*, and the obs. *abreyv*, represent it indirectly, through OFr. *abregier* and mid.Fr. *abre&vier*.

Like the latter, *abbreviate*, was often spelt *a-breviate* in 5-7.]  
 To make shorter, shorten, cut short in any way.

1530 PALSGR., I abreyvate: I make a thyng shorte, Je abrege. 1625 BACON *Essays* xxiv. 99 (1862) But it is one Thing to Abbreviate by Contracting, Another by Cutting off.

1 *trans.* To make a discourse shorter by omitting details and preserving the substance; to abridge, condense. *Obs.*

A. 1450 *Chester Pl.* i. 2 (Sh. Soc.) This mater he abbreviated into playes twenty-foure.

1592 GREENE *Conny catching* iii. 16 The queane abreviated her discourse. 1637 RALEIGH *Mahomet* 34 Abreviated out of two Arabique writers translated into Spanish. 1672 MANLEY *Interpreter* pref., I have omitted several Matters.contracted and abbreviated Others.

b To make an abstract or brief of, to epitomize. *Obs.*

C. 1450 *TREVISA Higden's Polychr.* i. 21 (Rolls Ser.) Trogus Pompeius, in hys xlii iii. bookes, allemoste of alle the storyes of the worlde, whom Iustinus his disciple did abreviate. 1603 FLORIO

*Montaigne* (1634) 627 To reade, to note, and to abbreviate Polibius. 1648-9 *The Kingdomes Weekly Intelligencer* Jan. 16 to 23 The high court of Justice did this day sit again concerning the trial of the King. The charge was brought in and abreviated.

c *Math.* To redu e (a fraction) to lower terms. *Obs.*

1796 *Mathem. Dict.* i. 2 To abbreviate fractions in arithmetic and algebra, is to lessen proportionally their terms, or the numerator and denominator.

2 *intr.* To speak or write briefly, to be brief. *Obs.*

1597 WARNER *Albion's Eng.* xli. lxxiv. 302 But new Rome left, of old Rome now abreuiat we will. 1622 MALYNES *Anc. Law-Merch.* 233 To abbreviate, I do referre the desirous Reader hereof to Master Hill his booke of Husbandrie.

3 *trans.* To shorten by cutting off a part; to cut short. Of time. *arch.*

1529 WHITINGTON *Vulgaria* 56 Ryot.abreviatedeth and shorteneth many a mannes lyfe. 1621 BURTON *Anat. Mel.* i. ii. 3. xv. 130 (1651) That adventure themselves and abbreviate their lives for the publike good. 1646 SIR T. BROWNE *Pseud. Ep.* 300 Against this we might very well set the length of their lives before the floud, which were abbreviated after.

b Of any operation occupying time.

1494 FAYAN VII. 333 If it sounde any thyng to theyr dishonoure, than shall it be abreuyatyd or hyd that the

abbreviate v.

**abbreviate**, v. Also 5-7  
 abbreviate. [f. abbreviate *ppl. a.*; or on the analogy of vbs. so formed; see -ate. A direct representative of L. *abbreviare*; as *abridge*, and the obs. *abreyv*, represent it indirectly, through OFr. *abregier* and mid.Fr. *abre&vier*. Like the latter, *abbreviate*, was often spelt *a-breviate* in 5-7.]

To make shorter, shorten, cut short in any way.

1 *trans.* To make a discourse shorter by omitting details and preserving the substance; to abridge, condense. *Obs.*

b To make an abstract or brief of, to epitomize. *Obs.*

c *Math.* To reduce (a fraction) to lower terms. *Obs.*

2 *intr.* To speak or write briefly, to be brief. *Obs.*

3 *trans.* To shorten by cutting off a part; to cut short. Of time. *arch.*

b Of any operation occupying time.

c Of things material; mostly *fig. arch.*

d Of words spoken or written, or symbols of any kind: To contract, so that a part stands for the whole. *The common mod. use.*

e Of sounds: To make (a vowel or syllable) short.

abbreviate v.

1530 PALSGR., I abreyvate: I make a thyng shorte, Je abrege.

1625 BACON *Essays* xxiv. 99 (1862) But it is one Thing to Abbreviate by Contracting, Another by Cutting off. A. 1450 *Chester Pl.* i. 2 (Sh. Soc.) This mater he abreviated into playes twenty-foure.

1592 GREENE *Conny catching* iii. 16 The queane abreviated her discourse.

1637 RALEIGH *Mahomet* 34 Abreviated out of two Arabique writers translated into Spanish.

1672 MANLEY *Interpreter* pref., I have omitted several Matters.contracted and abbreviated Others.

C. 1450 *TREVISA Higden's Polychr.* i. 21 (Rolls Ser.) Trogus Pompeius, in hys xlii iii. bookes, allemoste of alle the storyes of the worlde, whom Iustinus his disciple did abreviate.

1603 FLORIO *Montaigne* (1634) 627 To reade, to note, and to abbreviate Polibius.

1648-9 *The Kingdomes Weekly Intelligencer* Jan. 16 to 23 The high court of Justice did this day sit again concerning the trial of the King. The charge was brought in and abreviated.

1796 *Mathem. Dict.* i. 2 To abbreviate fractions in arithmetic and algebra, is to lessen proportionally their terms, or the numerator and denominator.

1597 WARNER *Albion's Eng.* xli. lxxiv. 302 But new Rome left, of old Rome now abreuiat we will.

1622 MALYNES *Anc. Law-Merch.* 233 To abbreviate, I do referre the desirous Reader hereof to Master Hill his booke of Husbandrie.

1529 WHITINGTON *Vulgaria* 56 Ryot.abreviatedeth and shorteneth many a mannes lyfe.

1621 BURTON *Anat. Mel.* i. ii. 3. xv. 130 (1651) That adventure themselves and abbreviate their lives for the publike good.

1646 SIR T. BROWNE *Pseud. Ep.* 300 Against this we might very well set the length of their lives before the floud, which were abbreviated after.

1494 FAYAN VII. 333 If it sounde any thyng to theyr dishonoure, than shall it be abreuyatyd or hyd that the trouthe shall not be known.

1655 FULLER *Ch. Hist.* ii. ix. 116 King Ethelbert was at his Devotions, which he would not omit, nor abbreviate for all their Clamour.

1665 E. B. TYLOR *Early Hist. Man.* iii. 46 The ancient Egyptian may be seen in the sculptures abbreviating the gesture.

1552 LATIMER *Serm. for 3rd Sund. in Adv. Wks.* ii. 287 His hand is not abreviated, or his power diminished.

1599 A. M. Gabelhouer's *Boock of Physicke* 178/2 Abbreviate es then the bagge, because it may genuelye, &

1661 MILTON *Accedence* (Wks. 1738) i. 607 The long way is much abbreviated, and the labour of understanding much more easy.

1508 SHAKS. *L.L.L.* v. 1. 26 He clepeth a Calf, Cause: Halfe, Haufe, neighbour vocatur nebour; neigh abreviated ne: this is abhominable.

1724 DE FOR etc. *A Tour* i. 364 (1769) The Exancoester of the Saxons, which was afterwards abbreviated to Excester and Exeter.

1888 GEIKIE *Phys. Geog.* i. iv. 27 Paris is situated two

abbreviate v.

**abbreviate**, v.

To make shorter, shorten, cut short in any way. 1530 Palsgr., I abreyvate: I make a thyng

1 *trans.*

b To make an abstract or brief of, to epitomize. *Obs.*C. 1450 *Trevisa Higden's*

c *Math.* To reduce (a fraction) to lower terms. *Obs.*1796 *Mathem. Dict.* i. 2 To

2 *intr.*

3 *trans.*

b Of any operation occupying time.1494 *Fayyan vii.* 333 If it sounde any thyng to

c Of things material; mostly *fig. arch.*1552 *Latimer Serm. for 3rd Sund. in Adv. Wks.* ii.

d Of words spoken or written, or symbols of any kind: To contract, so that a part

e Of sounds: To make (a vowel or syllable) short.1699 *Bentley Phalaris* 136 The

FIGURE 6. Dynamic Fragmentation of Abbreviate

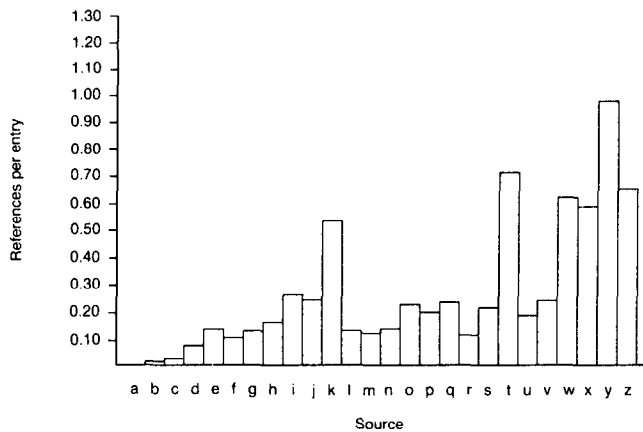


FIGURE 7a. Cross-references to Previous Letters

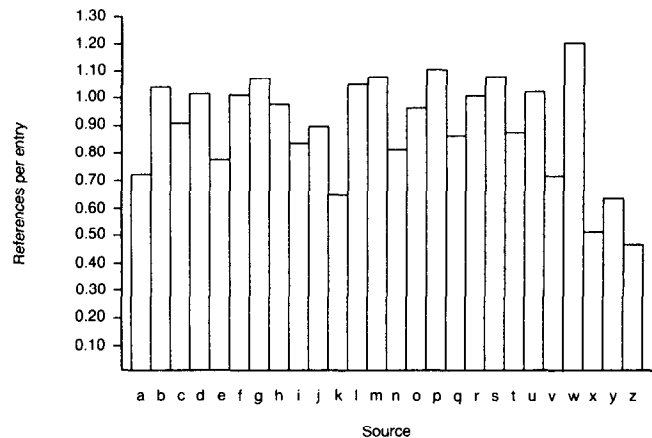


FIGURE 7b. Cross-references to Same Letter

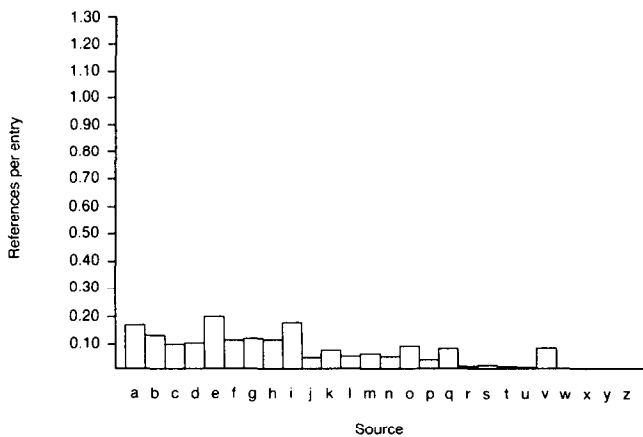


FIGURE 7c. Cross-references to Following Letters

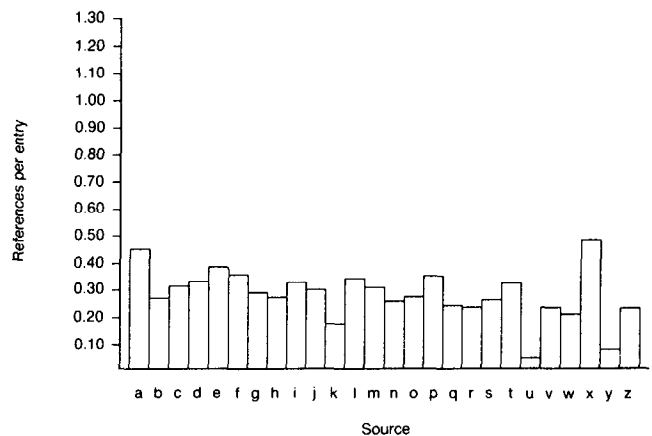


FIGURE 7d. Cross-references to Suffixes

should we link to *digging*, *out*, the phrase *digging out*, or perhaps all three? To provide for the third alternative is equivalent to isolating so-called *open compounds*, in general a difficult problem [1, 4]. Having determined the source, it may need to be reduced in tense or number to its morphological root. For example, *crevices* must be altered to *crevice* and *others* to *other*. On the other hand, because *digging* is explicitly defined in the *OED*, it should probably not be reduced to *dig*. Unlike most other dictionaries, the *OED* contains distinct entries for all derivative forms.

The third step is to determine the target, which could be a complete entry, a subentry for a derivative form, or perhaps a particular sense within an entry. In the above example, *washing-place* is listed as a derivative combination within the entry **washing** *vbl. sb.*, but only its subsense "(b) a place where gold is washed out from sand or earth" is directly applicable. Choosing the correct target involves determining the part of speech and etymological root to identify an entry and applying

sense disambiguation to identify a particular meaning [12]. For example, to identify the target of *pick*, one first has to determine that it is used as a noun rather than as a verb; next, that it is the first of four entries for **Pick** *sb.* that applies; and finally, identify which of the eight senses of this entry is the correct one.

The existing cross-reference links in the *OED* are important but limited in their scope. Furthermore, they can be treated as a special case of lexical links, which present several general difficulties. Because of the evolutionary development of the *OED*, not every cross-reference is specified precisely, so that resolving cross-references to words that appear in later volumes often involves the same target disambiguation required to solve lexical links in general. To support explicit intra-*OED* lexical links would therefore require sophisticated access software.

As in the case of text fragments, the tangle of problems involved in defining explicit links raises a question of their utility. In our experience, when rapid

searching and dynamic fragmentation of the resulting text is available, the user can usually determine the appropriate fragment by trying out a few alternatives. The cost of maintaining, updating, and storing explicit links would be considerable, and would not be offset by faster access than is provided by our present tools.

### RELATED PROBLEMS

In addition to evaluating hypertext for the *OED* itself, we have observed two related areas where hypertext could be used.

The first area is the use of the *OED* as a generator of hypertext links for other documents. The quotations in the *OED* can be interpreted as links from the dictionary to other texts. Links with a common source (i.e., those that start at the same sense of an entry) induce relationships among the referenced texts. Co-citation links could be derived for all texts that have some word sense in common; that is, texts are related if they each supply a quotation for a given sense. Useful links within and between other texts can also be identified from word collocations in the *OED*; for example, a large overlap among words in definitions can serve to identify similarity of topic [12]. Extracting meaningful pairings remains an interesting problem in computational linguistics.

A second area is the development and editing of entries for the *OED*. Computerized lexicography has a need for powerful hypertext-like structure editing, especially for the creation and maintenance of sense structure [15]. The source material for entries is a large set of quotations obtained over the years from volunteer and directed readers, each handwritten on a six-by-four-inch slip of paper. In a method still practiced today, the editors distill the senses of the word from the set of quotations by arranging and rearranging the slips into spatial categories on a desktop, looking for the pattern of historical development [14]. This approach seems highly susceptible to solution by a system that employs a desktop metaphor, as does NoteCards [9]. Existing systems, however, seem designed for tens of slips at a time, whereas *OED* editors can be faced with organizing thousands of slips. A more fundamental problem has been identified in experiments we have conducted on the organization of proverbs [17], which indicate that there is a quantifiable decrease in the quality of semantic categorization when a categorization task is performed in an electronic environment that employs a spatial metaphor. It appears that the ability to create temporary, unnamed categories is a key factor in the development of good semantic structures, and that current systems and metaphors interfere significantly with this ability.

### CONCLUSIONS

The investigation of cross-references in the *OED* showed significant local interconnectivity, but relatively few links between sections of the database that had been compiled at different times. We suspect this

localization of links to be a general tendency in large documents and possibly even in "hyper-libraries," simply because of the cost and difficulty of continually integrating old material with new.

Static fragmentation is common in hypertext systems, especially with data that has been created expressly for hypertext. Our experience with the *OED* indicates that static fragmentation is inappropriate for some converted texts. If there is too much implicit structure in the text or if the variance in structure is too great, then a better approach is to provide the user with tools for quickly developing dynamic fragmentations. Similarly, explicit storage and display of lexical links adds little to the user's ability to navigate the data, yet requires complex semantic analysis for proper resolution.

**Acknowledgments.** We gratefully acknowledge the financial assistance received from the Natural Sciences and Engineering Research Council of Canada through the University-Industry Program under grant 0039063.

### REFERENCES

1. Amsler, R. Research toward the development of a lexical knowledge base for natural language processing. Tech. Memo TM-ARH-010356. Bellcore, Morristown, N.J., Oct. 1987.
2. Benbow, T., Carrington, P., Johannesen, G., Tompa, F. W., and Weiner, E. Report on the *New Oxford English Dictionary* user survey. Tech. Rep. OED-87-05. UW Centre for the New Oxford English Dictionary, Waterloo, Ontario, Nov. 1987.
3. Berg, D.L., Gonnet, G.H., and Tompa, F.W. The New Oxford English Dictionary project at the University of Waterloo. In *Studies in Honour of Bernard Quemada*, A. Zampolli, Ed. Giardini Editore, Pisa, Italy, 1988.
4. Choueka, Y. Looking for needles in a haystack, or locating interesting collocational expressions in large textual databases. In *Proceedings of the RIAO (Recherche d'informations Assistée par Ordinateur) '88* (Cambridge, Mass., Mar. 21-24), 1988, pp. 609-623.
5. Conklin, J. Hypertext: An introduction and survey. *IEEE Computer* 20, 9 (Sept. 1987), 17-41.
6. Furnas, G.W. Generalized Fisheye Views. In *Proceedings of the CHI '86 Conference on Human Factors in Computing Systems* (Boston, Mass., Apr. 13-17). ACM, New York, 1986, pp. 16-23.
7. Gonnet, G.H., and Tompa, F.W. Mind your grammar: A new approach to modelling text. In *Proceedings of VLDB '87* (Brighton, England, Sept. 1-4), 1987, pp. 339-346.
8. Gonnet, G.H. Examples of PAT. Tech. Rep. OED-87-02. UW Centre for the New Oxford English Dictionary, Waterloo, Ontario, Aug. 1987.
9. Halasz, F.G., Moran, T.P., and Trigg, R.H. NoteCards in a nutshell. In *Proceedings of the CHI + GI Conference on Human Factors in Computer Systems and Graphics Interface* (Toronto, Ontario, Apr. 5-9). ACM, New York, 1987, pp. 45-52.
10. Jones, W.P. How do we distinguish the hyper from the hype in non-linear text?. In *Proceedings of INTERACT '87* (Stuttgart, Germany, Sept. 1-4). North-Holland, N. Y., 1987, pp. 1107-1113.
11. Kazman, R. Structuring the text of the Oxford English Dictionary through finite state transduction. Tech. Rep. CS-86-20. Dept. of Computer Science, University of Waterloo, Waterloo, Ontario, June 1986.
12. Lesk, M.E. Automatic sense disambiguation using machine readable dictionaries: How to tell a pine cone from an ice cream cone. In *Proceedings of the ACM SIGDOC Conference* (Toronto, Ontario, June 8-11). ACM, New York, 1986.
13. Murray, J.A.H., Bradley, H., Craigie, W.A., and Onions, C.T. *The Oxford English Dictionary*. Oxford at the Clarendon Press, Oxford, England, 1928.
14. Murray, K.M.E. *Caught in the Web of Words: James A.H. Murray and the Oxford English Dictionary*. Oxford University Press, Oxford, England, 1979.
15. Raymond, D.R., and Warburton, Y. Computerization of lexicographical activity on the New Oxford English Dictionary. Tech. Rep. OED-87-03. UW Centre for the New Oxford English Dictionary, Waterloo, Ontario, Dec. 10, 1986.



16. Raymond, D.R., and Blake, E.G. Solving queries in a grammar-defined *OED*. Unpublished tech. rep. UW Centre for the New Oxford English Dictionary, Waterloo, Ontario, Feb. 1987.
17. Raymond, D.R., Cañas, A.J., Tompa, F.W., and Safayeni, F.R. Measuring the effectiveness of personal database structures. *International J. Man-Machine Studies*. In press.
18. Weiner, E. The electronic English dictionary. *Oxford Magazine* (Feb. 1987), 6-9.

**CR Categories and Subject Descriptors:** H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing—*indexing methods*; *linguistic processing*; H.4 [Information Systems Applications]: Hypertext; J.5 [Arts and Humanities]: Dictionaries

**Additional Key Words and Phrases:** Document architecture, machine readable dictionaries

Authors' Present Addresses: Darrell R. Raymond and Frank Wm. Tompa, Dept. of Computer Science, University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

**Garg** (continued from p. 870)

17. Shaw, M. Abstraction techniques in modern programming languages. *IEEE Software* 1, 4 (Oct. 1984), 10-26.
18. Smith, J.M., and Smith, D.C.P. Database abstractions: Aggregations and generalizations. *ACM Trans. Database Systems*, 2, 2 (June 1977), 105-133.
19. Stoll, R.R. *Set Theory and Logic*. Dover, 1979.
20. Tichy, W. Design, implementation, and evaluation of a revision control system. In *6th International Conference on Software Engineering*, (Tokyo, Japan, 1982), IEEE Computer Society, 58-67.
21. Trigg, R.H. *A Network-Based Approach to Text Handling for the Online Scientific Community*. Ph.D. thesis, Maryland Artificial Intelligence Group, University of Maryland, November 1983.
22. *UNIX programmer's manual*. Bell Telephone Laboratories, Inc., Murray Hill, New Jersey.

**CR Categories and Subject Descriptors:** E.1 [Data Structures]: Hypertext; H.2.1 [Database Management]: Logical Design—*hypertext*; H.3.2 [Information Storage and Retrieval]: Information Storage—*hypertext*; I.7.2 [Text Processing]: Document Preparation

**General Terms:** Design, Theory

**Additional Key Words and Phrases:** Aggregation, generalization, and revision control

Author's Present Address: Pankaj K. Garg, Computer Science Department, University of Southern California, Los Angeles, CA 90089-0782.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

**Q:**

*What's Fast, Convenient, and Free?*




**A:**

ACM's "Order Express"  
Service for ACM Publications.

# 1-800-342-6626

Your credit card and our toll free number provide quick fulfillment of your orders.

- Journals
- Conference Proceedings
- SIG Newsletters
- SIGGRAPH VIDEO REVIEW
- "Computers in your Life" (An Introductory Film from ACM)

For Inquiries and other Customer Service call: (301) 528-4261

**acm** ASSOCIATION FOR COMPUTING MACHINERY